

## Accepted Manuscript

Deep metazoan phylogeny: When different genes tell different stories

Tetyana Nosenko, Fabian Schreiber, Maja Adamska, Marcin Adamski, Michael Eitel, Jörg Hammel, Manuel Maldonado, Werner E.G. Müller, Michael Nickel, Bernd Schierwater, Jean Vacelet, Matthias Wiens, Gert Wörheide

PII: S1055-7903(13)00029-8

DOI: <http://dx.doi.org/10.1016/j.ympev.2013.01.010>

Reference: YMPEV 4489

To appear in: *Molecular Phylogenetics and Evolution*

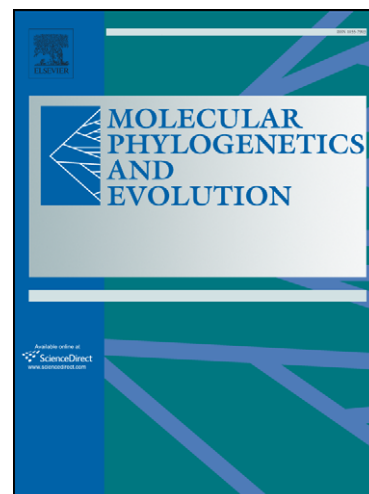
Received Date: 22 October 2012

Revised Date: 8 January 2013

Accepted Date: 12 January 2013

Please cite this article as: Nosenko, T., Schreiber, F., Adamska, M., Adamski, M., Eitel, M., Hammel, J., Maldonado, M., Müller, W.E.G., Nickel, M., Schierwater, B., Vacelet, J., Wiens, M., Wörheide, G., Deep metazoan phylogeny: When different genes tell different stories, *Molecular Phylogenetics and Evolution* (2013), doi: <http://dx.doi.org/10.1016/j.ympev.2013.01.010>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



**Deep metazoan phylogeny: When different genes tell different stories****Research Article**

Tetyana Nosenko<sup>a</sup>, Fabian Schreiber<sup>b</sup>, Maja Adamska<sup>c</sup>, Marcin Adamski<sup>c</sup>, Michael Eitel<sup>de</sup>, Jörg Hammel<sup>f</sup>, Manuel Maldonado<sup>g</sup>, Werner E.G. Müller<sup>h</sup>, Michael Nickel<sup>f</sup>, Bernd Schierwater<sup>e</sup>, Jean Vacelet<sup>i</sup>, Matthias Wiens<sup>h</sup>, and Gert Wörheide<sup>a,j,k</sup>

<sup>a</sup>Department of Earth and Environmental Sciences, Ludwig-Maximilians-Universität München, 80333 Munich, Germany

<sup>b</sup>Wellcome Trust Sanger Institute, Hinxton Hall, Hinxton, Cambridgeshire, CB10 1SA, UK

<sup>c</sup>Sars International Center for Marine Molecular Biology, 5008 Bergen, Norway

<sup>d</sup>Current address: Swire Institute of Marine Science, School of Biological Sciences, The University of Hong Kong, Hong Kong

<sup>e</sup>ITZ, Ecology and Evolution, Tierärztliche Hochschule Hannover, 30559 Hannover, Germany

<sup>f</sup>Institut of Systematic Zoology and Evolutionary Biology, Friedrich-Schiller-University of Jena, 07743 Jena, Germany

<sup>g</sup>Department of Marine Ecology, Centro de Estudios Avanzados de Blanes, 17300 Girona, Spain

<sup>h</sup>Institute of Physiological Chemistry University Medical Center, Johannes Gutenberg-University, 55128 Mainz, Germany

<sup>i</sup>CNRS UMR 7263 Institut Méditerranéen de Biodiversité et d'Ecologie Marine et continentale, Aix-Marseille Univ, 13007, Marseille, France

<sup>j</sup>GeoBio-Center<sup>LMU</sup>, Ludwig-Maximilians-Universität München, 80333 Munich, Germany

<sup>k</sup>Bayerische Staatssammlung für Paläontologie und Geologie, 80333 Munich,  
Germany

**Corresponding author:** Gert Wörheide, Department of Earth- and Environmental  
Sciences, Ludwig-Maximilians-University of Munich, Richard-Wagner-Str. 10,  
80333 Munich, Germany; Phone: +49 (0) 89 2180 6718; E-mail: woerheide@lmu.de

**Keywords:** Animal evolution; early metazoans; phylogeny; phylogenomics; gene  
selection

## Abstract

Molecular phylogenetic analyses have produced a plethora of controversial hypotheses regarding the patterns of diversification of non-bilaterian animals. To unravel the causes for the patterns of extreme inconsistencies at the base of the metazoan tree of life, we constructed a novel supermatrix containing 122 genes, enriched with non-bilaterian taxa. Comparative analyses of this supermatrix and its two non-overlapping multi-gene partitions (including ribosomal and non-ribosomal genes) revealed conflicting phylogenetic signals. We show that the levels of saturation and Long Branch Attraction artifacts in the two partitions correlate with gene sampling. The ribosomal gene partition exhibits significantly lower saturation levels than the non-ribosomal one. Additional systematic errors derive from significant variations in amino acid substitution patterns among the metazoan lineages that violate the stationarity assumption of evolutionary models frequently used to reconstruct phylogenies. By modifying gene sampling and the taxonomic composition of the outgroup, we were able to construct three different yet well-supported phylogenies. These results show that the accuracy of phylogenetic inference may be substantially improved by selecting genes that evolve slowly across the Metazoa and applying more realistic substitution models. Additional sequence-independent genomic markers are also necessary to assess the validity of the phylogenetic hypotheses.

## 1. Introduction

The historical sequence of early animal diversification events has been the subject of debate for approximately a century. Morphological character analyses leave a degree of uncertainty concerning the evolutionary relationships among the five major

metazoan lineages: Porifera, Placozoa, Ctenophora, Cnidaria, and Bilateria (Collins et al., 2005). In the last few years, this debate has been fueled by a plethora of conflicting phylogenetic hypotheses generated using molecular data (Dunn et al., 2008; Erwin et al., 2011; Philippe et al., 2009; Pick et al., 2010; Schierwater et al., 2009; Sperling et al., 2009). The persisting controversy includes questions concerning the earliest diverging animal lineage (Porifera vs. Placozoa vs. Ctenophora), the validity of the Eumetazoa (Bilateria + Cnidaria + Ctenophora) and Coelenterata (Cnidaria + Ctenophora) clades, and relationships among the main lineages of Porifera (sponges; reviewed in Wörheide et al., 2012). These questions are fundamental for understanding the evolution of both animal body plans and genomes (Philippe et al., 2009).

In 2003, Rokas and co-authors (Rokas et al., 2003a) showed that the evolutionary relationships between major metazoan lineages cannot be resolved using single genes or a small number of protein-coding sequences. Because of the high stochastic error, the analyses of the individual genes resulted in conflicting phylogenies. These authors also observed that at least 8,000 randomly selected characters (> 20 genes) are required to overcome the effect of these discrepancies (Rokas et al., 2003b).

However, the authors' subsequent attempt at resolving the deep metazoan relationships using a large dataset containing 50 genes from 17 metazoan taxa (including six non-bilaterian species) was not successful (Rokas et al., 2005). By contrast, the analysis of the identical set of genes robustly resolved the higher-level phylogeny of Fungi, a group of approximately the same age as the Metazoa (Yuan et al., 2005). Based on this result, these authors concluded that because of the rapidity of the metazoan radiation, the true phylogenetic signal preserved on the deep internal branches was too low to reliably deduce their branching order (Rokas and Carroll,

2006). However, this conclusion did not discourage scientists from further attempts at resolving this difficult phylogenetic question using the traditional sequence-based phylogenetic approach. The main strategy of the subsequent studies was increasing the amount of data, including both gene and taxon sampling. In 2008, a novel hypothesis of early metazoan evolution was proposed by Dunn et al. (2008) based on the analysis of 150 nuclear genes (21,152 amino acid [aa] characters) from 71 metazoan taxa (however, with only nine non-bilaterian species among them). According to this hypothesis, ctenophores represent the most ancient, earliest diverging branch of the Metazoa. This evolutionary scenario did not gain any support from the analysis of another large alignment that contained 128 genes (30,257 aa) and a larger number of non-bilateral metazoan species (22; Philippe et al., 2009). This study revived the Coelenterata and Eumetazoa hypotheses (Hyman, 1940) and placed the Placozoa as a sister group of the Eumetazoa. Another scenario for early metazoan evolution was proposed by Schierwater et al. (2009) based on the analysis of a dataset that included not only nuclear protein-coding genes but also mitochondrial genes and morphological characters (a “total evidence” dataset). This study reconstructed monophyletic “Diploblasta” (*i.e.*, non-bilaterian metazoans) with a “basal” Placozoa as a sister group of the Bilateria.

Recently published metazoan phylogenies differ in their taxon and gene sampling and their application of phylogenetic methods and thresholds, including the use of different models of amino acid substitution. Any of these factors may be a source of the observed incongruity among the proposed deep metazoan phylogenies (Dunn et al., 2008; Philippe et al., 2009; Schierwater et al., 2009). Comparative analyses of the three above-described multi-gene alignments showed that the observed conflict can be partially attributed to the presence of contaminations, alignment errors, and reliance

on simplified evolutionary models (Philippe et al., 2011) or long branch attraction artifacts caused by insufficient ingroup taxon sampling (Pick et al., 2010). Correcting the alignment errors in the datasets by Dunn *et al.* (2008) and Schierwater et al. (2009) and applying an evolutionary model that best fit these data, altered both the tree topology and basal node support, but failed to resolve the incongruences between the three phylogenies.

The objective of the present study is to further assess the causes of inconsistency between deep (non-bilaterian) metazoan phylogenies obtained using phylogenomic (large multi-gene) datasets with a main emphasis on the effect of gene sampling. We approached this question with multiple comparative analyses of a novel phylogenomic dataset with two multi-gene sub-matrices that have identical taxon samplings, comparable lengths, and missing data percentage but different gene contents. We also increased the taxon sampling by adding new data from non-bilaterian lineages, including seven Porifera species, one Ctenophora species, and a novel placozoan strain.

## 2. Materials and Methods

### 2.1. Data acquisition

New data were generated for nine species of non-bilaterian metazoans, including one ctenophore, *Beroe* sp., an unidentified placozoan species (Placozoan sp. strain H4), and seven sponges: *Asbestopluma hypogea*, *Ephydatia muelleri*, *Pachydictyum globosum*, *Tethya wilhelma* (all from class Demospongiae), *Crateromorpha meyeri* (class Hexactinellida), *Corticium candelabrum* (class Homoscleromorpha), (Expressed Sequence Tag [EST] libraries), and *Sycon ciliatum* (class Calcarea; EST

and genomic data). The data generation information and complete list of taxa included in the analyses are provided in Supplementary Materials.

## 2.2. Multi-gene matrix assembly

A total of 225 orthologous groups (OGs) dominated by non-bilaterian taxa were constructed using the automated ortholog assignment pipeline OrthoSelect (Schreiber et al., 2009). The input data used by OrthoSelect consisted of complete genome and EST data for 71 species, including 21 species of Porifera, two placozoans, four ctenophores, 13 cnidarians, 21 bilaterians, three choanoflagellates, two ichthyosporeans, one filasterean, and four species of Fungi (Supplementary Dataset S1). The OGs containing less than 40 taxa were discarded from the analysis. Due to an uneven distribution of complete genome sequence data among the species included in our dataset, these OGs were dominated by sequences for bilaterian and outgroup taxa. To minimize the effect of alignment construction artifacts (*e.g.*, misalignments, paralogous and contaminant sequences) on phylogenetic inference, the remaining OGs were further processed using the following three-step procedure:

*step I.* Paralog and contamination pruning. Sequences in each OG were aligned using the computer program MUSCLE v3.8 (Edgar, 2004) and annotated using a sequence similarity search (BLAST; *e*-value threshold  $10^{-20}$ ) against the NCBI nr. Paralogous and contaminant sequences were identified and removed from the OGs based on the result of the BLAST annotation and a visual inspection of the motives conserved among all taxa in the alignment. After this procedure, all OGs containing less than 40 taxa were discarded from the analysis. The remaining OGs were re-aligned with MUSCLE. Ambiguously aligned regions were removed with TrimAl v1.2 (Capella-Gutiérrez et al., 2009) using a heuristic selection of the trimming method based on similarity statistics. This program allows for a coordinated trimming of multiple



alignments according to the consistency score inferred from the most conserved alignments. The resulting alignments were refined manually (*e.g.*, by correcting small frameshifts and removing the remaining ambiguously aligned sites);

*step II.* Identifying paralogous and contaminant sequences in each OG using a tree-based approach modified from Rodriguez-Ezpeleta et al. (2007). Briefly, each OG was analyzed under the CAT+ $\Gamma$ 4 model using PhyloBayes version 3.2e (Lartillot et al., 2007; Lartillot and Philippe, 2004, 2006). Markov Chain Monte Carlo (MCMCs) were run for 11,000 cycles. Posterior consensus trees were constructed for each gene after discarding the initial 3000 cycles. The sequences that formed well-supported sub-clusters that conflicted with both super-matrix trees, produced long branches, or were “trapped” by a distant outgroup (Filasterea, Ichthyosporea, or Fungi) were excluded from individual gene alignments as paralogous or contaminant. The OGs containing less than 40 taxa were excluded from further analyses; and

*step III.* The compositional homogeneity test implemented in PhyloBayes was conducted for each OG using chains obtained during the step II. All OGs that did not pass the compositional deviation score threshold ( $z < 2$ ) were discarded (see Supplementary Dataset S2).

After the OG cleaning and filtering, the most distant outgroup, Fungi, which served as a trap for the contaminant sequences, was excluded from the alignments to reduce the computing time and LBA artifact.

The 122 OGs that passed the three-step selection procedure (Supplementary Dataset S2) were classified by function according to the KOG database functional classification (Tatusov et al., 2003) and sorted into two groups. One group included 87 genes encoding proteins involved in translation (ribosomal proteins). We emphasize that ribosomal RNA genes, which have frequently been used for

reconstructing metazoan phylogenies (Mallatt et al., 2012; Medina et al., 2001; Peterson and Eernisse, 2001), were not included in this dataset. The remaining 35 OGs from different functional classes formed the second dataset hereafter termed non-ribosomal. The single-gene OGs were concatenated using FASconCAT (Kuck and Meusemann, 2010) to obtain the 14615 aa-long ribosomal, 9187 aa-long non-ribosomal, and 22975 aa-long combined multi-gene matrices (Table 1 and Supplementary Dataset S2). To reduce the ribosomal-to-non-ribosomal site ratio in the second combined dataset, 2731 ribosomal sites (nine genes) represented by less than 38 ingroup taxa were removed from the alignment (20244 aa-long combined multi-gene matrix; Supplementary Dataset S1).

### *2.3. Taxon sampling and missing data*

The resulting datasets were used to construct several sub-matrices (Table 1) that differed by taxon sampling size (42-67 taxa) and percentage of missing data. The datasets were constructed under three different missing-data-per-taxon thresholds: 50, 80, and 95%. The total amount of missing characters varied from 14 to 36% across datasets. The largest ribosomal and non-ribosomal datasets (Table 1) were constructed under the relaxed missing data cutoff stringency, in which up to 95% missing data were allowed per taxon for lineages represented by more than two species. After the exclusion of all outgroup taxa but choanoflagellates, the dataset consisted of 63 taxa. To test the effect of taxon sampling (and missing data) on the tree topology and basal node support, we excluded the following taxa from the 14615 aa-long ribosomal dataset: (I) seven bilaterian species containing higher amounts of missing data (2-3 from each major bilaterian lineage; 56-taxa matrix); (II) all species containing more than 50% missing data (49-taxa matrix); and (III) all species containing more than 50% missing data and the same seven bilaterian species as in matrix I (42-taxa matrix;

see Supplementary Dataset S1).

To reduce the missing data effect and computing time, the seven bilaterian species and all non-bilaterian taxa containing more than 80% missing data were excluded from all 50-taxa matrices (ribosomal, non-ribosomal, and combined; Supplementary Dataset S1) used for phylogenetic analyses. The missing data threshold used in this study was established at 30% total characters (Table 1). The only dataset that had a higher percentage of missing data (36%), the 63-taxa non-ribosomal gene matrix, was used solely for assessing the taxon sampling and missing data effects.

#### *2.4. Evolutionary model selection*

The choice of model of protein evolution is well-known to affect the pattern of phylogenetic relationships among major metazoan lineages inferred from molecular data (Jeffroy et al., 2006; Philippe et al., 2011). To select the model that best fit our data, we analyzed each of the 122 OGs using ProtTest (Abascal et al., 2005). The fit of the LG model for the concatenated ribosomal and non-ribosomal matrices compared to more complex evolutionary models, which are not available under the Maximum Likelihood framework (GTR, CAT, and CAT-GTR), was assessed using a cross-validation test (Stone, 1974). The cross-validation test was conducted using PhyloBayes as described in Supplementary Materials.

#### *2.5. Phylogenetic Analyses*

ML trees were obtained with RAxML v7.2.7 (Stamatakis et al., 2005) under the LG model (Le and Gascuel, 2008). Bayesian analyses were performed using PhyloBayes v3.2e and the CAT, CAT-GTR, LG, and GTR models. The taxon-specific compositional heterogeneities were estimated under the CAT model using the algorithm implemented in PhyloBayes. The patristic- and p-distances for the saturation analyses were computed using PATRISTIC (Fourment and Gibbs, 2006)

and MEGA5 (Tamura et al., 2011), respectively. To identify taxa that have the most unstable phylogenetic position in our trees, we conducted leaf stability analyses (Thorley and Wilkinson, 1999) using Phyutility (Smith and Dunn, 2008). The full details and descriptions of the techniques above are provided in Supplementary Materials.

The new sequence data reported in this paper were deposited in GenBank (<http://www.ncbi.nlm.nih.gov>) and the European Nucleotide Archive (ENA; <http://www.ebi.ac.uk/ena>); the alignments were deposited at OpenDataLMU (<http://dx.doi.org/10.5282/ubm/data.55>).

### 3. Results

#### 3.1. *Different gene matrices tell different stories*

The ProtTest analyses indicated that LG+ $\Gamma$ +I was the evolutionary model that best fit the majority of the single-gene alignments in a Maximum Likelihood (ML) framework. However, a further statistical comparison (cross-validation test; Stone, 1974) extended to more complex evolutionary models rejected the LG in favor of GTR (scores of 383 and 61 in favor of GTR for the ribosomal and non-ribosomal matrices, respectively), which, in turn, was outperformed by both the Bayesian CAT (with a score difference of 1027 for the ribosomal and 1219 for non-ribosomal matrices) and CAT-GTR (1239 and 1264) models. Although CAT-GTR was identified as the best model for these data, most of our analyses were conducted using the CAT model because of computational constraints. To illustrate the problem, 20,000 cycles of MCMCs run for our ribosomal gene matrix containing 63 taxa and

14,615 aa positions were completed in 48 days under the CAT model, whereas runs under the CAT-GTR model required 202 days to complete.

The phylogenetic analyses of the most data-rich supermatrix, which contains 122 genes (22,975 sites) and 50 taxa (Table 1), under the CAT model is presented in Fig. 1. We used the sister group of the Metazoa in this analysis, the Choanoflagellata (King et al., 2008), as the only outgroup. This tree supports the Coelenterata and monophyly of sponges but provides no resolution for the relationships between Coelenterata, Porifera, and Bilateria. In addition, the placement of Placozoa as the sister group of the Porifera is not well supported. The lack of resolution for the deep nodes in this tree reflects major conflicts between the previously published metazoan phylogenies (Dunn et al., 2008; Philippe et al., 2009; Pick et al., 2010; Schierwater et al., 2009; Sperling et al., 2009). To identify the source of the potential conflict within this dataset, we divided this matrix into two non-overlapping multi-gene partitions (Supplementary Dataset S2). One partition included 87 genes (14,615 sites) from a single functional class: translation (primarily ribosomal proteins). Another partition consisted of 35 genes (9,187 sites) that represented 11 functional classes. The phylogenetic analyses of the two partitions resulted in incongruent topologies (Figs. 2A, B, and 3). The analyses of the ribosomal gene matrices under the CAT model output a well-resolved tree that provided strong support for the Coelenterata and Eumetazoa concepts and monophyly of Porifera (Fig. 2A). The only basal node that did not receive high support was the Placozoa and Porifera divergence. The analysis of the ribosomal datasets conducted under the CAT-GTR model was consistent with that conducted under the CAT model on phylum-level relationships, including the monophyly of Porifera (Supplementary Fig. S1). In addition, this analysis provided strong support for Placozoa as the sister group of the Porifera. However, the best-

fitting model left the relative positions of the Bilateria, Coelenterata, and Placozoa-Porifera clades unresolved. No apparent misplacement of taxa (including those containing over 80% missing data) was observed in these phylogenies. Reducing the taxon sampling by selectively excluding species from only bilaterian clades, only non-bilaterian clades, or both, did not alter the tree topologies but led to a gradual decrease in the support values at the deep nodes under both the CAT and CAT-GTR models (Figs. 2A and S1).

Unlike the ribosomal trees, the topology of the non-ribosomal tree rooted with choanoflagellates was sensitive to missing data. The Bayesian analysis of the non-ribosomal gene matrix containing 63 taxa under the CAT model resulted in several misplacements of taxa containing more than 80% missing data and, consequently, poor support for the phylum-level nodes (*e.g.*, Bilateria and Cnidaria; Supplementary Fig. S2). Therefore, the “gappy” taxa were removed from the non-ribosomal and combined matrices. The topology of the non-ribosomal tree containing 50 taxa was not consistent with the ribosomal CAT and CAT-GTR trees on the relationships of the deep branches. This topology disrupts the monophyly of sponge lineages, does not support Coelenterata, and determines Ctenophora to be the sister group to the remaining Metazoa. We emphasize that this “Ctenophora-basal” topology was common in all of our rooted non-ribosomal ML and Bayesian trees constructed under the LG, GTR, CAT-GTR, and CAT models.

To further assess the effect of gene sampling on the higher-level metazoan phylogeny, we decreased the proportion of ribosomal sites by excluding nine ribosomal genes (12% of the combined matrix length) from the combined dataset. The resulting matrix contained 54% ribosomal and 46% non-ribosomal sites. This modification restored Coelenterata and its sister relationships with Bilateria (99% PP)

but broke the Porifera-Placozoa group into three paraphyletic clades: Placozoa, Calcarea-Homoscleromorpha, and Demospongiae-Hexactinellida (Fig. 4). The Placozoa were recovered as the sister group to the Eumetazoa. Unlike the original tree depicted in Fig. 1, all basal nodes of this “shortened matrix” tree received strong PP support ( $\geq 95\%$ ).

### 3.2. Saturation and Long Branch Attraction (LBA) artifacts

Saturation and LBA are two factors that may contribute to the instability of the metazoan phylogeny observed in this study and explain its sensitivity to gene sampling (Bergsten, 2005; Philippe et al., 2011; Pick et al., 2010). We conducted the following tests to assess whether the above-described conflicts in tree topology (*e.g.*, the position of the Ctenophora and relationships among the Porifera lineages) resulted from saturation and LBA: (I) measured the relative saturation levels in the ribosomal and non-ribosomal partitions; (II) analyzed a less saturated matrix under the models of protein evolution that fit these data less well than the CAT model; (III) removed all non-metazoan taxa from the two datasets and constructed un-rooted trees under the CAT model; and (IV) replaced the Choanoflagellata with a more distant outgroup and reconstructed the ribosomal and non-ribosomal phylogenies under the CAT model.

To compare the saturation levels in our ribosomal and non-ribosomal gene matrices, we plotted the patristic distances inferred from the corresponding trees against the uncorrected p-distances (Fig. 3). The results of this test revealed a higher saturation level in the non-ribosomal gene matrix (the regression line slope = 0.36 and Pearson correlation coefficient  $R = 0.26$ ) compared to our ribosomal gene dataset (slope = 0.42;  $R = 0.84$ ; an ideal non-saturated dataset has a slope = 1 and  $R = 1$ ).

We next assumed that if the topology inferred from the non-ribosomal gene matrix under the CAT model resulted from saturation, it should be reproducible with

a less saturated matrix and less well-fitting model. To test this prediction, we analyzed our ribosomal gene matrix using two standard evolutionary models: the LG and GTR. These models have been shown to be more susceptible to saturation and LBA artifacts (Lartillot and Philippe, 2004) and fit our data less well than the CAT model. The outcome was consistent with our prediction: the “Ctenophora-basal” and paraphyletic Porifera were recovered in all ribosomal trees constructed under the LG and GTR models (Supplementary Fig. S3). This result strongly suggests that a similar position of these branches in the non-ribosomal CAT trees is likely to be an artifact of a higher saturation level in this gene set, which increases the branch length variance and potentially adds to an LBA bias (Felsenstein, 1978).

To test for an LBA bias, we excluded the non-metazoan outgroup taxa from the analysis as the most obvious source of LBA (Holland et al., 2003) and constructed unrooted ribosomal, non-ribosomal, and “combined” CAT trees. The removal of the choanoflagellates resolved most conflicts between the resulting phylogenies. In all three un-rooted phylogenies, the ctenophores and cnidarians tended to establish sister group relationships, with weaker support from the non-ribosomal dataset, however. Regarding the sponges, the Silicea *sensu stricto* (Demospongiae + Hexactinellida) represent the sister-group to the Homoscleromorpha + Calcarea clade (Supplementary Fig. S4). Obviously, the issue of sponge mono- vs. paraphyly depends on where the root of the tree is placed.

Another standard method for detecting LBA artifacts is to use distant outgroups (reviewed in Bergsten, 2005). A distant outgroup increases the LBA effect and works as a trap for the long ingroup branches. Previous analyses by Philippe *et al.* (2009) demonstrated that including the additional outgroups distantly related to Metazoa (in particular, Filasterea, Ichthyosporea, and Fungi) into their dataset reduced the support



values for the deep metazoan nodes. We used a slightly different approach to identify the ingroup branches affected by LBA. Instead of increasing the outgroup size, we replaced the choanoflagellates with Ichthyosporea, a group of organisms more distant from the Metazoa than the Choanoflagellata and Filasterea (Shalchian-Tabrizi et al., 2008; Torruella et al., 2012). This replacement led to major rearrangements in both the ribosomal and non-ribosomal trees (Supplementary Fig. S5A, B). The position of the ctenophores in the non-ribosomal tree did not change. Instead, this branch switched to the base of the Metazoa in the less saturated ribosomal tree. In addition, the Cnidaria-Bilateria clade was disrupted in both phylogenies. Now, both Coelenterata lineages appeared at the basal position to other animals in the non-ribosomal tree and, presumably as a consequence of this shift, the monophyly of Porifera and its sister group relationships with the Placozoa were restored with a high level of support (Supplementary Fig. S5B).

The results of these tests demonstrate a strong effect of LBA by the outgroup on metazoan tree topology, including inter- and intra-phyla level relationships. The extent of this effect depends on the saturation level in the given multi-gene matrix (as determined by gene sampling), choice of outgroup, and assumptions of the evolutionary model used in the analysis.

### 3.3. Leaf stability and among-taxa compositional heterogeneity

One of the methods commonly applied to diminish systematic error and biases is to exclude unstable taxa and those that have a biochemical composition significantly deviating from the global empirical composition of the dataset (Brinkmann and Philippe, 1999; Thorley and Wilkinson, 1999). To identify taxa that have an unstable phylogenetic position in our ribosomal and non-ribosomal trees, we calculated leaf stability (LS) indices (Thorley and Page, 2000) for all species using the Bayesian

CAT trees sampled during the MCMC chains. According to the results of the LS analysis, all representatives of Homoscleromorpha, Calcarea, and Placozoa were unstable in all of our trees. Choanoflagellates, ichthyosporeans, filastereans, and ctenophores received low LS values from several datasets (Supplementary Table S1). In addition, the posterior predictive analysis of among-taxa compositional heterogeneity showed that the amino acid composition of the choanoflagellate, ichthyosporean, filastereans, and placozoan sequences deviated significantly from the global empirical biochemical composition in both datasets (Supplementary Table S1). Potentially, the presence of the above-mentioned taxa in the alignments increases LBA and destabilizes the resulting phylogeny. The analyses of the LBA artifacts presented above confirmed a destabilizing effect of choanoflagellates and ichthyosporeans on metazoan trees. High (relative to metazoans) alanine and low lysine contents in both outgroup taxa and high glycine and low leucine contents in ichthyosporeans indicate that compositional heterogeneity can be partially attributed to high GC content in both outgroups (King et al., 2008; Codon Usage Database; Supplementary Fig. S5C and D). However, excluding the placozoans, the most unstable ingroup lineage (Supplementary Table S1), from the analysis changed neither the topology of the non-ribosomal tree, nor that of the ribosomal tree (data not shown).

#### **4. Discussion**

##### *4.1. Why do different genes tell different stories?*

The multiple conflicting metazoan phylogenies presented here and in previous publications (Dunn et al., 2008; Erwin et al., 2011; Philippe et al., 2009; Pick et al., 2010; Schierwater et al., 2009; Sperling et al., 2009; Srivastava et al., 2010) have one

feature in common: they have long terminal and short internal branches. Frequently, such a topology is a sign of ancient rapid radiations, which are closely spaced diversification events that occurred deep in time (Rokas et al., 2003a; Rokas et al., 2005). This observation is consistent with both the fossil record and molecular clock estimates showing that the radiation of early metazoans occurred within a relatively short time span of approximately 700 MYA (Erwin et al., 2011). A major challenge of phylogenetic reconstructions associated with such ancient and likely rapid radiations is recovering the true signal at the deep nodes. Previously published studies showed that sequence alignments containing one or few genes provide information insufficient for resolving the relationships between major metazoan lineages (Rokas et al., 2003a). Our results are consistent with this conclusion: none of the 122 single-gene alignments constructed for this study provide any support for the deep nodes. Increasing the size of the dataset (both taxon and gene sampling) has been thought to be the logical solution since at least 8000 randomly selected characters are required to obtain reasonable support for ancient diversifications (Rokas et al., 2003b). Owing to recent advances in DNA sequencing technologies, considerable amounts of sequence data are available for constructing phylogenomic alignments consisting of hundreds of genes. However, there is an uncertainty regarding the best gene sampling strategy. A common practice is the *a posteriori* sampling of as many genes shared by the lineages of interest as the data allow (Dunn et al., 2008; Gatesy and Baker, 2005; Kuck and Meusemann, 2010; Srivastava et al., 2010). This method minimizes heuristic and other cognitive biases associated with *a priori* choice of target genes. However, the method is based on the assumption that the collective phylogenetic signal from all OGs should be stronger than noise (Hillis, 1998). This assumption is often violated when phylogenetic problems associated with ancient rapid radiations

are addressed (Bergsten, 2005). The analysis of different partitions of a phylogenomic alignment is the most reliable method to assess the validity of this assumption for a particular dataset. The consistency of phylogenies inferred from independent partitions remains the strongest evidence of an accuracy of phylogenetic estimates (Comas et al., 2007; Swofford, 1991).

In this study, we used the partitioning of a large alignment to test the effect of gene sampling on the higher-level metazoan phylogeny and assess the validity of the random-gene sampling strategy in application to this problem. There are several possible approaches for defining multi-gene partitions, such as gene-specific evolutionary rates, linkage, and gene function (Miyamoto and Fitch, 1995). Partitioning based on evolutionary rates is a promising approach that would test the prediction that slow evolving genes are the most suitable for resolving ancient diversifications, whereas more rapidly evolving genes should be selected for testing recent radiation events (Donoghue and Sanderson, 1992; Felsenstein, 1983; Giribet, 2002). In phylogenomics, relative evolutionary rates are estimated either based on single-gene saturation plots or by calculating the length of each gene tree (the sum of all branch lengths) or pairwise sequence distances (Bevan et al., 2005; Ebersberger et al., 2011; Fong and Fujita, 2011; Graybeal, 1994). However, these methods are not reliable when comparing single-gene alignments containing different amounts of missing data. Since complete genome sequences are available for few non-bilaterian metazoan species, the alignments used in this study (and in other genomic-scale deep metazoan phylogeny studies) are dominated by EST-derived sequences and contain relatively high amounts of missing data (13-36% missing data in our matrices and 50 and 27% in the datasets from Dunn et al. (2008) and Philippe et al. (2009), respectively (Table 1). In this study, we partitioned our total dataset based on gene

functions as a proxy for the rate of evolution (reviewed in Koonin and Wolf, 2006). We constructed two non-overlapping matrices sufficiently long for analyzing deep metazoan phylogeny (> 8,000 characters; as suggested by Rokas et al., 2003b). One matrix exclusively included the housekeeping genes involved in translation, which are highly conserved and show uniformly slow rates of evolution across the Tree of Life (Castillo-Davis et al., 2004; Hori et al., 1977; Hughes et al., 2006; Landais et al., 2003; Moreira et al., 2002; Warren et al., 2010). Because of the ubiquitously high expression levels, these genes can be found in EST libraries of all organisms and therefore constitute a significant component of phylogenomic alignments constructed to address higher-level metazoan phylogeny (*e.g.*, 26 and 11% of all sites in the supermatrices by Dunn *et al.*, 2008, and Philippe *et al.*, 2009, respectively). The second partition was constructed in accordance with the “randomness” criterion. This partition included genes from various functional categories characterized by various rates of evolution from slow evolving ubiquitins and histones (an evolutionary rate similar to ribosomal proteins) to less constrained metabolic enzymes (Nei et al., 2000; Piontkivska et al., 2002; Rooney et al., 2002). The phylogenetic analyses of the two partitions produced conflicting trees (Fig. 2). Moreover, combining the genes from the two datasets in different proportions either led to a loss of the basal-node support (Fig. 1) or resulted in a well-supported topology incongruent with the two partition trees (Fig. 4). This surprisingly high sensitivity of the non-bilaterian component of the metazoan phylogeny to gene sampling may result from different levels of non-phylogenetic signal in our datasets. Since all gene alignments were constructed using the same methods and selected using the same statistical tests and thresholds (described in Materials and Methods), all matrices were expected to have similar levels of systematic error associated with ortholog selection and aligning. The results

of saturation and LBA tests indicate that these artifacts provide the most plausible explanation for the observed inconsistency of the resulting phylogenies. The dataset that included genes from various functional categories had a significantly higher saturation level than the ribosomal-gene matrix (Fig. 4). The phylogenies generated using this “random-gene” matrix exhibited stronger LBA biases (*e.g.*, the basal position of the Ctenophora relative to other metazoan lineages in all rooted trees) than the phylogenies generated using the ribosomal gene dataset. This result is consistent with the prediction that limiting analyses to slow evolving genes is the best approach for resolving high-level phylogenies. Depending on the history and rate of evolution, genes are known to vary in their phylogenetic informativeness over historical time (Felsenstein, 1983; Graybeal, 1994). Sites informative for resolving the relationships between the terminal branches can be homoplasious at deeper nodes of a phylogenetic tree. Restraining the analyses to genes that evolve slowly across the Tree of Life may reduce the level of saturation in the dataset and recover the phylogenetic signal at the basal nodes. This conclusion does not contradict and instead complements the “randomness” criterion. However, this conclusion assumes a significant reduction of the number of candidate genes and consequently, restrains the character sampling (length) of the deep metazoan phylogenomic datasets.

Although our ribosomal tree depicted in Fig. 2A received high statistical support for the basal nodes and showed no apparent LBA effect and a low sensitivity to taxon sampling, the distant outgroup test and CAT-GTR analysis revealed a degree of relationship instability among the Bilateria, Coelenterata, and Placozoa-Porifera branches (Supplementary Fig. S1). This instability can be attributed to low-level biases in the ribosomal trees. Saturation and LBA biases result from the substantial variation of evolutionary processes both along a sequence and among the lineages

(Lartillot and Philippe, 2004; Lopez et al., 2002). Problems occur when this variation violates the assumptions of the evolutionary model used. Although genes that have the most heterogeneous biochemical composition were excluded from our datasets, the comparison of the taxon-specific amino acid frequencies revealed a significant among-lineage compositional deviation in both partitions. In particular, ctenophores, placozoans, and outgroup taxa exhibited biochemical compositions that significantly deviated from the global empirical amino acid frequencies in both alignments (Supplementary Table S1). The factors that contribute to among-lineage compositional heterogeneity include a historical shift in site-specific substitution rates and qualitative changes of substitution patterns over time (Lopez et al., 2002; Roure and Philippe, 2011). The models used in this study (and the other studies on higher-level metazoan phylogeny cited above) account for the across-site heterogeneity but assume a homogeneous evolutionary process over time (Lartillot and Philippe, 2004). The patterns observed in both of our datasets violate this assumption and provide an additional source of systematic error, which may contribute to the observed instability of the early metazoan phylogeny.

To summarize, this study generated three incongruent, yet strongly supported tree topologies: the ribosomal gene tree (Fig. 2A), the combined dataset II tree (Fig. 4), and the non-ribosomal gene tree containing an ichthyosporean outgroup (Supplementary Fig. S5B). The latter phylogeny can be rejected with high confidence because it was based on the most saturated dataset and was not confirmed by the analysis with the outgroup closest to the Metazoa. The remaining two datasets have their advantages and disadvantages. The combined dataset is longer than the ribosomal one and includes genes from various functional categories and is therefore less prone to gene sampling bias. However, the level of saturation in this dataset is

increased due to the inclusion of the non-ribosomal matrix. The ribosomal gene matrix has the lowest saturation level. The resulting phylogeny is robust to the alterations of taxon sampling. The main criticism of the ribosomal gene phylogeny is that it is based on functionally coupled macromolecules, which might share a common evolutionary bias (Bleidorn et al., 2009). Apparently this tree reflects the early evolution of translational machinery in animals. The question is whether the history of the metazoan translation machinery is congruent with its species phylogeny. Answering this question is particularly important for resolving the position of the Placozoa and the relationships between the major sponge lineages.

Although our phylogenetic reconstructions left a degree of uncertainty regarding the relationships among the early branching animal clades, the dynamics of the tree topology changes under the different models and with different outgroups shed light on several controversies of the metazoan phylogeny.

#### *4.2. Phylogenetic positions of the Placozoa and Porifera lineages*

Recently published hypotheses on the phylogenetic position of the placozoans include but are not limited to the Placozoa basal to other eumetazoans (Philippe et al., 2009; Srivastava et al., 2008; Srivastava et al., 2010), the Placozoa as a sister group of Bilateria (Pick et al., 2010), Bilateria-Cnidaria (Ryan et al., 2010), or Coelenterata-Porifera clades (Schierwater et al., 2009). The relationships among the major Porifera lineages represent another point of conflict among the metazoan trees. Several studies indicate sponges as a paraphyletic group (Dunn et al., 2008; Erwin et al., 2011; Medina et al., 2001; Peterson and Eernisse, 2001; Rokas et al., 2005; Sperling et al., 2009); other studies argue for the monophyly of Porifera (Philippe et al., 2009, 2011; Pick et al., 2010; reviewed in Wörheide *et al.*, 2012). The phylogenetic patterns observed in this study link these two phylogenetic problems together. All of our trees



supporting sponge monophyly place Placozoa as a sister group of Porifera (Figs. 2A, 3, S1, and S5A, B), whereas the paraphyletic sponges always coincide with placozoans placed as a sister-group of eumetazoans (Figs. 2B and S3). Our less-saturated dataset analyzed under the best-fitting models favors the first scenario (sponge monophyly; Figs. 2A and S1). However, regardless of the tree topology and confidence values for the corresponding nodes, the phylogenetic positions of Placozoa, Homoscleromorpha, and Calcarea are extremely unstable (Supplementary Table S1). In addition to a significantly deviating amino acid composition and a global interplay among the long and short branches of the tree, the factors that may contribute to the observed instability include an uneven distribution of taxon sampling (Bergsten, 2005; Hillis, 1998). Although we added new taxa to all lineages listed above, these groups apparently remain undersampled. Our taxon sampling test shows that support for the Porifera monophyly increases when the taxon sampling increases (Fig. 2A). Based on this observation, we predict that adding new species of calcareous sponges and homoscleromorphs should increase the stability of the Porifera clade and potentially resolve its relationships with Placozoa.

#### *4.3. Ctenophora as the most problematic branch among the non-bilaterians*

Morphological and molecular studies gave rise to several controversial hypotheses on the phylogenetic position of ctenophores (Dunn et al., 2008; Wallberg et al., 2004). In this study, we obtained trees supporting two hypotheses: the ctenophores as a sister group of Cnidaria (Coelenterata hypothesis, Figs. 1, 2A, and 4; Haeckel, 1866) and the ctenophores basal to other animals (“Ctenophora-early” hypothesis, Figs. 2B, S3, and S5; Dunn et al., 2008). The comparison of our ribosomal and non-ribosomal gene phylogenies generated under different models of evolution provides several supporting arguments that the position of ctenophores as a sister-group to the

remaining Metazoa in our trees is an artifact of LBA between the outgroup and ctenophore branches: (I) Ctenophora-basal did not receive strong support in any tree analyzed under the CAT model when the Choanoflagellata, the closest to the Metazoa lineage, was used as an outgroup. This position of ctenophores was supported either when the trees were generated under a less-fitting amino acid substitution model or a more distant outgroup was used (Supplementary Figs. S3 and S5); and (II) in the absence of non-metazoan taxa, the unrooted ribosomal and non-ribosomal phylogenies were consistent with the sister group relationships between Ctenophora and Cnidaria (Supplementary Fig. S4).

The ctenophores consistently formed long branches in all Bayesian and ML trees constructed for this study. Poor ctenophore taxon sampling may partially explain the problem. Large sequence datasets (EST libraries) are available for only four ctenophore species. Including taxa that represent the overall diversity of a problematic group in the phylogenetic datasets is perceived as the most efficient method of breaking up long branches (Hillis, 1998). However, the lack of a robust ctenophore taxonomy (Podar et al., 2001) and insufficient knowledge of their biology (in particular, the rates of self-fertilization in hermaphroditic ctenophores) challenge the development of an efficient taxon sampling strategy. Self-fertilization is associated with high mutation rates (Schultz and Lynch, 1997); therefore, the presence of self-fertilized species in phylogenetic datasets may increase saturation and aggravate the LBA problem (Pett et al., 2011).

Another concern is that the long branch separating the ctenophores from their closest living relatives may indicate an extensive extinction of ancient ctenophore forms. The hypothesis that all extant ctenophore species evolved from a relatively recent common ancestor was proposed by Podar et al. (2001) based on the

phylogenetic analyses of 18S sequences from 26 ctenophore species. This assumption is also supported by the fossil record, in which putative stem-group Ctenophores from the Cambrian differ from recent taxa in a number of manners (*e.g.*, the number of comb rows, presence of lobate organs in the former, *etc.*) and represent extinct stem groups (Carlton et al., 2007; King et al., 2008). Our results do not contradict this hypothesis. The evolutionary distances between four species, each representing one of the major ctenophore lineages, are short in comparison to those between the major lineages of sponges and cnidarians (Figs. 1, 2, and 4). If this hypothesis is true, ctenophores may be the most problematic branch of the non-bilaterian section of the metazoan tree and be difficult to resolve even with additional taxon sampling.

## 5. Conclusions

This study shows an extreme sensitivity of the higher-level metazoan phylogeny to the gene composition of the phylogenomic matrices. The gene sampling strategy determines the level of saturation and LBA biases in the resulting phylogenies. According to our results, a careful *a priori* (*i.e.*, post-sequencing and before analyses) selection of genes that evolve slowly across all metazoan lineages helps to decrease systematic errors and recover the phylogenetic signal from the noise. Using this approach, we were able to reconstruct a metazoan phylogeny that is consistent with traditional, morphology-based views on the phylogeny of non-bilaterian metazoans, including monophyletic Porifera and ctenophores as a sister group of cnidarians. The stability of the metazoan tree can be further improved by applying a more realistic amino acid substitution model that accounts for the variation of evolutionary rates and biochemical patterns, both along the sequences and among the lineages, and by increasing the taxon sampling of critically “undersampled” lineages. In the case of

non-bilaterian animals, these lineages should be drawn from calcareous and homoscleromorph sponges, placozoans, and ctenophores. In addition, identifying and sampling early branching, slowly evolving outgroup species with an amino acid composition similar to the metazoan ingroup may help to decrease the outgroup effect.

The above steps promise to significantly improve the robustness of deep phylogeny estimation. However, the criteria used to assess the fit and performance of new evolutionary models and validity of the resulting phylogeny remain to be identified. In this study, we confirmed the previous conclusion that the standard measures of clade support, such as the bootstrap or Bayesian posterior probabilities, may support several conflicting hypotheses with high apparent confidence. When different multi-gene partitions tell different stories, we cannot rely solely on traditional phylogenetic analyses of long (and even longer) sequences. Difficult phylogenetic problems, such as the relationships between the major metazoan lineages, call for the development of new, sequence-independent genomic markers (SIGMs, *e.g.*, protein domain architecture, gene order, gene fusions, duplications, insertions-deletions, or genetic code variants; Rokas and Holland, 2000) that would provide independent data to test conflicting phylogenetic hypotheses. Although attempts to use such markers, for example microRNAs to resolve sponge relationships (Sperling et al., 2010), transposable elements (short interspersed elements, SINEs; Piskurek and Jackson, 2011) and changes in spliceosomal intron positions (NIPs; Lehmann et al., 2012), to resolve early metazoan relationships have thus far been unsuccessful, the growing number of fully sequenced genomes of non-bilaterian animals might provide sufficient data in the future to discover novel SIGMs to test phylogenomic hypotheses and finally enable us to fully appreciate the early evolution of animals.

### **Acknowledgments**

We thank S. Leys, B. Bergum, Ch. Arnold, M. Krüß, and E. Gaidos for providing samples; M. Kube and his team (MPE for Molecular Genetics, Berlin, Germany) for library construction; I. Ebersberger and his team (Center for Integrative Bioinformatics, Vienna, Austria) for data processing; and K. Nosenko for the artwork. This work was financially supported by the German Research Foundation (DFG Priority Program SPP1174 “Deep Metazoan Phylogeny,” projects Wo896/6 and WI 2216/2-2). M.A. and Mn.A. acknowledge funding from Sars International Centre for Marine Molecular Biology and the Research Council of Norway. M.E. acknowledges financial support by the Evangelisches Studienwerk e.V. Villigst and the German Academic Exchange service (DAAD).

### **Author contributions**

G.W. conceived the research and obtained the funding; T.N. and G.W. designed the research; T.N. and F.S. analyzed the data; M.A., Mn.A., M.E., J.H., B.S., W.M., M.W. and G.W. provided data; M.M., M.N., and J.V. provided samples; M.M. contributed to manuscript revision; and T.N. and G.W. wrote the paper.

### **Appendices A, B, and C. Supplementary Materials**

Supplementary materials associated with this article can be found in the online version, at (*the link*).

### **References**

Abascal, F., Zardoya, R., Posada, D., 2005. ProfTest: selection of best-fit models of protein evolution. *Bioinformatics* 21, 2104-2105.

- Bergsten, J., 2005. A review of long-branch attraction. *Cladistics* 21, 163-193.
- Bevan, R.B., Lang, B.F., Bryant, D., 2005. Calculating the evolutionary rates of different genes: a fast, accurate estimator with applications to maximum likelihood phylogenetic analysis. *Syst Biol* 54, 900-915.
- Bleidorn, C., Podsiadlowski, L., Zhong, M., Eeckhaut, I., Hartmann, S., Halanych, K.M., Tiedemann, R., 2009. On the phylogenetic position of Myzostomida: can 77 genes get it wrong? *BMC Evol Biol* 9, 150.
- Brinkmann, H., Philippe, H., 1999. Archaea sister group of Bacteria? Indications from tree reconstruction artifacts in ancient phylogenies. *Mol Biol Evol* 16, 817-825.
- Capella-Gutiérrez, S., Silla-Martínez, J., Gabaldón, T., 2009. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25, 1972-1973.
- Carlton, J.M., Hirt, R.P., Silva, J.C., Delcher, A.L., Schatz, M., Zhao, Q., Wortman, J.R., Bidwell, S.L., Alsmark, U.C., Besteiro, S., et al., 2007. Draft genome sequence of the sexually transmitted pathogen *Trichomonas vaginalis*. *Science* 315, 207-212.
- Castillo-Davis, C.I., Kondrashov, F.A., Hartl, D.L., Kulathinal, R.J., 2004. The functional genomic distribution of protein divergence in two animal phyla: coevolution, genomic conflict, and constraint. *Genome Res* 14, 802-811.
- Collins, A.G., Cartwright, P., McFadden, C.S., Schierwater, B., 2005. Phylogenetic context and Basal metazoan model systems. *Integr Comp Biol* 45, 585-594.
- Comas, I., Moya, A., Gonzalez-Candelas, F., 2007. From phylogenetics to phylogenomics: the evolutionary relationships of insect endosymbiotic gamma-Proteobacteria as a test case. *Syst Biol* 56, 1-16.
- Donoghue, M., Sanderson, M., 1992. The suitability of molecular and morphological evidence in reconstructing plant phylogeny. In: Soltis, P., Soltis, D., Doyle, J. (Eds.), *Molecular Systematics in Plants*. Chapman and Hall, New York, pp. 340-368.
- Dunn, C.W., Hejnal, A., Matus, D.Q., Pang, K., Browne, W.E., Smith, S.A., Seaver, E., Rouse, G.W., Obst, M., Edgecombe, G.D., et al., 2008. Broad phylogenomic sampling improves resolution of the animal tree of life. *Nature* 452, 745-749.
- Ebersberger, I., de Matos Simoes, R., Kupczok, A., Gube, M., Kothe, E., Voigt, K., von Haeseler, A., 2011. A consistent phylogenetic backbone for the fungi. *Mol Biol Evol* 29, 1319-1334.
- Edgar, R.C., 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32, 1792-1797.
- Erwin, D.H., Laflamme, M., Tweedt, S.M., Sperling, E.A., Pisani, D., Peterson, K.J., 2011. The Cambrian conundrum: early divergence and later ecological success in the early history of animals. *Science* 334, 1091-1097.
- Felsenstein, J., 1978. A likelihood approach to character weighting and what it tells us about parsimony and compatibility. *Biol J Linn Soc* 16, 183-196.
- Felsenstein, J., 1983. Parsimony in systematics: Biological and statistical issues. *Annu Rev Ecol Evol S* 14, 313-333.
- Fong, J.J., Fujita, M.K., 2011. Evaluating phylogenetic informativeness and data-type usage for new protein-coding genes across Vertebrata. *Mol Phylogenet Evol* 61, 300-307.
- Fourment, M., Gibbs, M.J., 2006. PATRISTIC: a program for calculating patristic distances and graphically comparing the components of genetic change. *BMC Evol Biol* 6, 1.

- Gatesy, J., Baker, R.H., 2005. Hidden likelihood support in genomic data: can forty-five wrongs make a right? *Syst Biol* 54, 483-492.
- Giribet, G., 2002. Current advances in the phylogenetic reconstruction of metazoan evolution. A new paradigm for the Cambrian Explosion? *Molecular Phylogenetics and Evolution* 24, 345-357.
- Graybeal, A., 1994. Evaluating the Phylogenetic Utility of Genes - a Search for Genes Informative About Deep Divergences among Vertebrates. *Syst Biol* 43, 174-193.
- Haeckel, E., 1866. *Generelle Morphologie der Organismen*. G. Reimer Berlin.
- Hillis, D.M., 1998. Taxonomic sampling, phylogenetic accuracy, and investigator bias. *Syst Biol* 47, 3-8.
- Holland, B.R., Penny, D., Hendy, M.D., 2003. Outgroup misplacement and phylogenetic inaccuracy under a molecular clock - A simulation study. *Syst Biol* 52, 229-238.
- Hori, H., Higo, K., Osawa, S., 1977. The rates of evolution in some ribosomal components. *J Mol Evol* 9, 191-201.
- Hughes, J., Longhorn, S.J., Papadopoulou, A., Theodorides, K., de Riva, A., Mejia-Chang, M., Foster, P.G., Vogler, A.P., 2006. Dense taxonomic EST sampling and its applications for molecular systematics of the Coleoptera (beetles). *Mol Biol Evol* 23, 268-278.
- Hyman, L., 1940. *The Invertebrates: Protozoa through Ctenophora*. McGraw-Hill, New York.
- Jeffroy, O., Brinkmann, H., Delsuc, F., Philippe, H., 2006. Phylogenomics: the beginning of incongruence? *Trends Genet* 22, 225-231.
- King, N., Westbrook, M.J., Young, S.L., Kuo, A., Abedin, M., Chapman, J., Fairclough, S., Hellsten, U., Isogai, Y., Letunic, I., et al., 2008. The genome of the choanoflagellate *Monosiga brevicollis* and the origin of metazoans. *Nature* 451, 783-788.
- Koonin, E.V., Wolf, Y.I., 2006. Evolutionary systems biology: links between gene evolution and function. *Curr Opin Biotechnol* 17, 481-487.
- Kuck, P., Meusemann, K., 2010. FASconCAT: Convenient handling of data matrices. *Mol Phylogenet Evol* 56, 1115-1118.
- Landais, I., Ogliaastro, M., Mita, K., Nohata, J., Lopez-Ferber, M., Duonor-Cerutti, M., Shimada, T., Fournier, P., Devauchelle, G., 2003. Annotation pattern of ESTs from *Spodoptera frugiperda* Sf9 cells and analysis of the ribosomal protein genes reveal insect-specific features and unexpectedly low codon usage bias. *Bioinformatics* 19, 2343-2350.
- Lartillot, N., Brinkmann, H., Philippe, H., 2007. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. *BMC Evol Biol* 7 Suppl 1, S4.
- Lartillot, N., Philippe, H., 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol Biol Evol* 21, 1095-1109.
- Lartillot, N., Philippe, H., 2006. Computing Bayes factors using thermodynamic integration. *Syst Biol* 55, 195-207.
- Le, S.Q., Gascuel, O., 2008. An improved general amino acid replacement matrix. *Mol Biol Evol* 25, 1307-1320.
- Lehmann, J., Stadler, P.F., Krauss, V., 2012. Near intron pairs and the metazoan tree. *Mol Phylogenet Evol.*, doi: 10.1016/j.ympev.2012.11.012
- Lopez, P., Casane, D., Philippe, H., 2002. Heterotachy, an important process of protein evolution. *Mol Biol Evol* 19, 1-7.

- Mallatt, J., Waggoner-Craig, C., Yoder, M.J., 2012. Nearly complete rRNA genes from 371 Animalia: Updated structure-based alignment and detailed phylogenetic analysis. *Mol Phylogenet Evol* 64, 603-617.
- Medina, M., Collins, A.G., Silberman, J.D., Sogin, M.L., 2001. Evaluating hypotheses of basal animal phylogeny using complete sequences of large and small subunit rRNA. *Proc Natl Acad Sci U S A* 98, 9707-9712.
- Miyamoto, M.M., Fitch, W.M., 1995. Testing species phylogenies and phylogenetic methods with congruence. *Syst. Biol.* 44, 64-76.
- Moreira, D., Kervestin, S., Jean-Jean, O., Philippe, H., 2002. Evolution of eukaryotic translation elongation and termination factors: Variations of evolutionary rate and genetic code deviations. *Mol Biol Evol* 19, 189-200.
- Nei, M., Rogozin, I.B., Piontkivska, H., 2000. Purifying selection and birth-and-death evolution in the ubiquitin gene family. *Proc Natl Acad Sci U S A* 97, 10866-10871.
- Peterson, K.J., Eernisse, D.J., 2001. Animal phylogeny and the ancestry of bilaterians: inferences from morphology and 18S rDNA gene sequences. *Evol Dev* 3, 170-205.
- Pett, W., Ryan, J.F., Pang, K., Mullikin, J.C., Martindale, M.Q., Baxevanis, A.D., Lavrov, D.V., 2011. Extreme mitochondrial evolution in the ctenophore *Mnemiopsis leidyi*: Insight from mtDNA and the nuclear genome. *Mitochondrial DNA* 22, 130-142.
- Philippe, H., Brinkmann, H., Lavrov, D.V., Littlewood, D.T., Manuel, M., Worheide, G., Baurain, D., 2011. Resolving difficult phylogenetic questions: why more sequences are not enough. *PLoS Biol* 9, e1000602.
- Philippe, H., Derelle, R., Lopez, P., Pick, K., Borchiellini, C., Boury-Esnault, N., Vacelet, J., Renard, E., Houliston, E., Queinnec, E., et al., 2009. Phylogenomics revives traditional views on deep animal relationships. *Curr Biol* 19, 706-712.
- Pick, K.S., Philippe, H., Schreiber, F., Erpenbeck, D., Jackson, D.J., Wrede, P., Wiens, M., Alie, A., Morgenstern, B., Manuel, M., et al., 2010. Improved phylogenomic taxon sampling noticeably affects nonbilaterian relationships. *Mol Biol Evol* 27, 1983-1987.
- Piontkivska, H., Rooney, A.P., Nei, M., 2002. Purifying selection and birth-and-death evolution in the histone H4 gene family. *Mol Biol Evol* 19, 689-697.
- Piskurek, O., Jackson, D.J., 2011. Tracking the ancestry of a deeply conserved eumetazoan SINE domain. *Mol Biol Evol* 28, 2727-2730.
- Podar, M., Haddock, S.H., Sogin, M.L., Harbison, G.R., 2001. A molecular phylogenetic framework for the phylum Ctenophora using 18S rRNA genes. *Mol Phylogenet Evol* 21, 218-230.
- Rodriguez-Ezpeleta, N., Brinkmann, H., Burger, G., Roger, A.J., Gray, M.W., Philippe, H., Lang, B.F., 2007. Toward resolving the eukaryotic tree: the phylogenetic positions of jakobids and cercozoans. *Curr Biol* 17, 1420-1425.
- Rokas, A., Carroll, S.B., 2006. Bushes in the Tree of Life. *PLoS Biology* 4: e352.
- Rokas, A., Holland, P.W., 2000. Rare genomic changes as a tool for phylogenetics. *Trends Ecol Evol* 15, 454-459.
- Rokas, A., King, N., Finnerty, J., Carroll, S.B., 2003a. Conflicting phylogenetic signals at the base of the metazoan tree. *Evol Dev* 5, 346-359.
- Rokas, A., Kruger, D., Carroll, S.B., 2005. Animal evolution and the molecular signature of radiations compressed in time. *Science* 310, 1933-1938.
- Rokas, A., Williams, B.L., King, N., Carroll, S.B., 2003b. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature* 425, 798-804.



- Rooney, A.P., Piontkivska, H., Nei, M., 2002. Molecular evolution of the nontandemly repeated genes of the histone 3 multigene family. *Mol Biol Evol* 19, 68-75.
- Roure, B., Philippe, H., 2011. Site-specific time heterogeneity of the substitution process and its impact on phylogenetic inference. *BMC Evol Biol* 11, 17.
- Ryan, J.F., Pang, K., Mullikin, J.C., Martindale, M.Q., Baxevanis, A.D., 2010. The homeodomain complement of the ctenophore *Mnemiopsis leidyi* suggests that Ctenophora and Porifera diverged prior to the ParaHoxozoa. *EvoDevo* 1, 9.
- Schierwater, B., Eitel, M., Jakob, W., Osigus, H.J., Hadrys, H., Dellaporta, S.L., Kolokotronis, S.O., Desalle, R., 2009. Concatenated analysis sheds light on early metazoan evolution and fuels a modern "urmetazoon" hypothesis. *PLoS Biol* 7, e20.
- Schreiber, F., Pick, K., Erpenbeck, D., Worheide, G., Morgenstern, B., 2009. OrthoSelect: a protocol for selecting orthologous groups in phylogenomics. *BMC Bioinformatics* 10, 219.
- Schultz, S.T., Lynch, M., 1997. Deleterious mutation and extinction: effects of variable mutational effects, synergistic epistasis, beneficial mutations, and degree of outcrossing. *Evolution* 51, 1363-1371.
- Shalchian-Tabrizi, K., Minge, M.A., Espelund, M., Orr, R., Ruden, T., Jakobsen, K.S., Cavalier-Smith, T., 2008. Multigene phylogeny of choanozoa and the origin of animals. *PLoS One* 3, e2098.
- Smith, S.A., Dunn, C.W., 2008. Phyutility: a phyloinformatics tool for trees, alignments and molecular data. *Bioinformatics* 24, 715-716.
- Sperling, E.A., Peterson, K.J., Pisani, D., 2009. Phylogenetic-signal dissection of nuclear housekeeping genes supports the paraphyly of sponges and the monophyly of Eumetazoa. *Mol Biol Evol* 26, 2261-2274.
- Sperling, E.A., Robinson, J.M., Pisani, D., Peterson, K.J., 2010. Where's the glass? Biomarkers, molecular clocks, and microRNAs suggest a 200-Myr missing Precambrian fossil record of siliceous sponge spicules. *Geobiology* 8, 24-36.
- Srivastava, M., Begovic, E., Chapman, J., Putnam, N.H., Hellsten, U., Kawashima, T., Kuo, A., Mitros, T., Salamov, A., Carpenter, M.L., et al., 2008. The Trichoplax genome and the nature of placozoans. *Nature* 454, 955-960.
- Srivastava, M., Simakov, O., Chapman, J., Fahey, B., Gauthier, M.E., Mitros, T., Richards, G.S., Conaco, C., Dacre, M., Hellsten, U., et al., 2010. The Amphimedon queenslandica genome and the evolution of animal complexity. *Nature* 466, 720-726.
- Stamatakis, A., Ludwig, T., Meier, H., 2005. RAxML-III: a fast program for maximum likelihood-based inference of large phylogenetic trees. *Bioinformatics* 21, 456-463.
- Stone, M., 1974. Cross-validators choice and assessment of statistical prediction. *J R Stat Soc Ser B* 36, 111-147.
- Swofford, D., 1991. When are phylogeny estimates from molecular and morphological data incongruent? In: Miyamoto, M.M., Cracraft, J. (Eds.), *Phylogenetic Analysis of DNA Sequences*. Oxford University Press, pp. 295-333.
- Tamura, K., Peterson, D., Peterson, N., Stecher, G., Nei, M., Kumar, S., 2011. MEGA5: Molecular Evolutionary Genetics Analysis Using Maximum Likelihood, Evolutionary Distance, and Maximum Parsimony Methods. *Mol Biol Evol* 28, 2731-2739.
- Tatusov, R.L., Fedorova, N.D., Jackson, J.D., Jacobs, A.R., Kiryutin, B., Koonin, E.V., Krylov, D.M., Mazumder, R., Mekhedov, S.L., Nikolskaya, A.N., et al.,

2003. The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* 4, 41.
- Thorley, J.L., Page, R.D., 2000. RadCon: phylogenetic tree comparison and consensus. *Bioinformatics* 16, 486-487.
- Thorley, J.L., Wilkinson, M., 1999. Testing the phylogenetic stability of early tetrapods. *J Theor Biol* 200, 343-344.
- Torruella, G., Derelle, R., Paps, J., Lang, B.F., Roger, A.J., Shalchian-Tabrizi, K., Ruiz-Trillo, I., 2012. Phylogenetic relationships within the Opisthokonta based on phylogenomic analyses of conserved single-copy protein domains. *Mol Biol Evol* 29, 531-544.
- Wallberg, A., Thollesson, M., Farris, J.S., Jondelius, U., 2004. The phylogenetic position of the comb jellies (Ctenophora) and the importance of taxonomic sampling. *Cladistics* 20, 558-578.
- Warren, A.S., Anandakrishnan, R., Zhang, L., 2010. Functional bias in molecular evolution rate of *Arabidopsis thaliana*. *BMC Evol Biol* 10, 125.
- Wörheide, G., Dohrmann, M., Erpenbeck, D., Larroux, C., Maldonado, M., Voigt, O., Borchellini, C., Lavrov, D.V., 2012. Deep Phylogeny and Evolution of Sponges (Phylum Porifera). In: Becerro, M.A., Uriz, M.J., Maldonado, M., Turon, X. (Eds.), *Advances in Marine Biology*. Vol. 61, Academic Press, The Netherlands: Amsterdam, pp. 1-78.
- Yuan, X., Xiao, S., Taylor, T.N., 2005. Lichen-like symbiosis 600 million years ago. *Science* 308, 1017-1020.

Codon Usage Database. <http://www.kazusa.or.jp/codon/>.

## Tables

Table 1. Large multi-gene matrices used for addressing the early metazoan phylogeny question.

## Figure legends

Fig. 1. Bayesian consensus tree inferred from the analysis of the matrix composed of both ribosomal and non-ribosomal genes (22975 aa positions and 50 terminal taxa) under the CAT+ $\Gamma$  model. The solid circles indicate nodes that received maximum Posterior Probabilities support (PP 100%). Numbers are given for nodes that have PP < 100% (PP < 95% is given in italics). The scale bar indicates the number of changes per site.

Fig. 2. Comparative analyses of two multi-gene partitions. **(A)** Bayesian consensus tree inferred from the analysis of the ribosomal gene partition containing 14615 aa positions and 63 terminal taxa. The PPs were obtained from the analyses of the ribosomal sub-matrices containing 63, 56, 49, and 42 taxa (Table 1). The solid circles indicate maximum PP support (100%) from all datasets. The blue color indicates species excluded from the 56- and 42-taxa sub-matrices; the red color indicates species excluded from the 49- and 42-taxa sub-matrices. Due to the conflicting relative positions of *mertensiid* sp. 3 and *Pleurobrachia pileus* in different trees, the corresponding node was collapsed. **(B)** Bayesian consensus tree inferred from the analysis of the non-ribosomal gene partition containing 9187 amino acid positions and 50 terminal taxa. The PP and scale bars are as in Fig. 1. All trees were constructed under the CAT+ $\Gamma$  model.

Fig. 3. Saturation analysis. The relative saturation levels were estimated for the ribosomal and non-ribosomal gene matrices containing 50 taxa by computing the Pearson correlation coefficient  $R$  and slope of the regression line of patristic vs.  $p$ -distances. The patristic distances between pairs of taxa were inferred from the branch lengths of ML trees constructed under the LG+ $\Gamma$ 8+I model.

Fig. 4. Bayesian consensus trees obtained from the analyses of the combined matrix II (20244 aa positions and 50 taxa; Table 1) under the CAT+ $\Gamma$  model. This matrix differs from the combined matrix I (Fig. 1) by 2731 ribosomal sites. The PP and scale bar are as in Fig. 1.

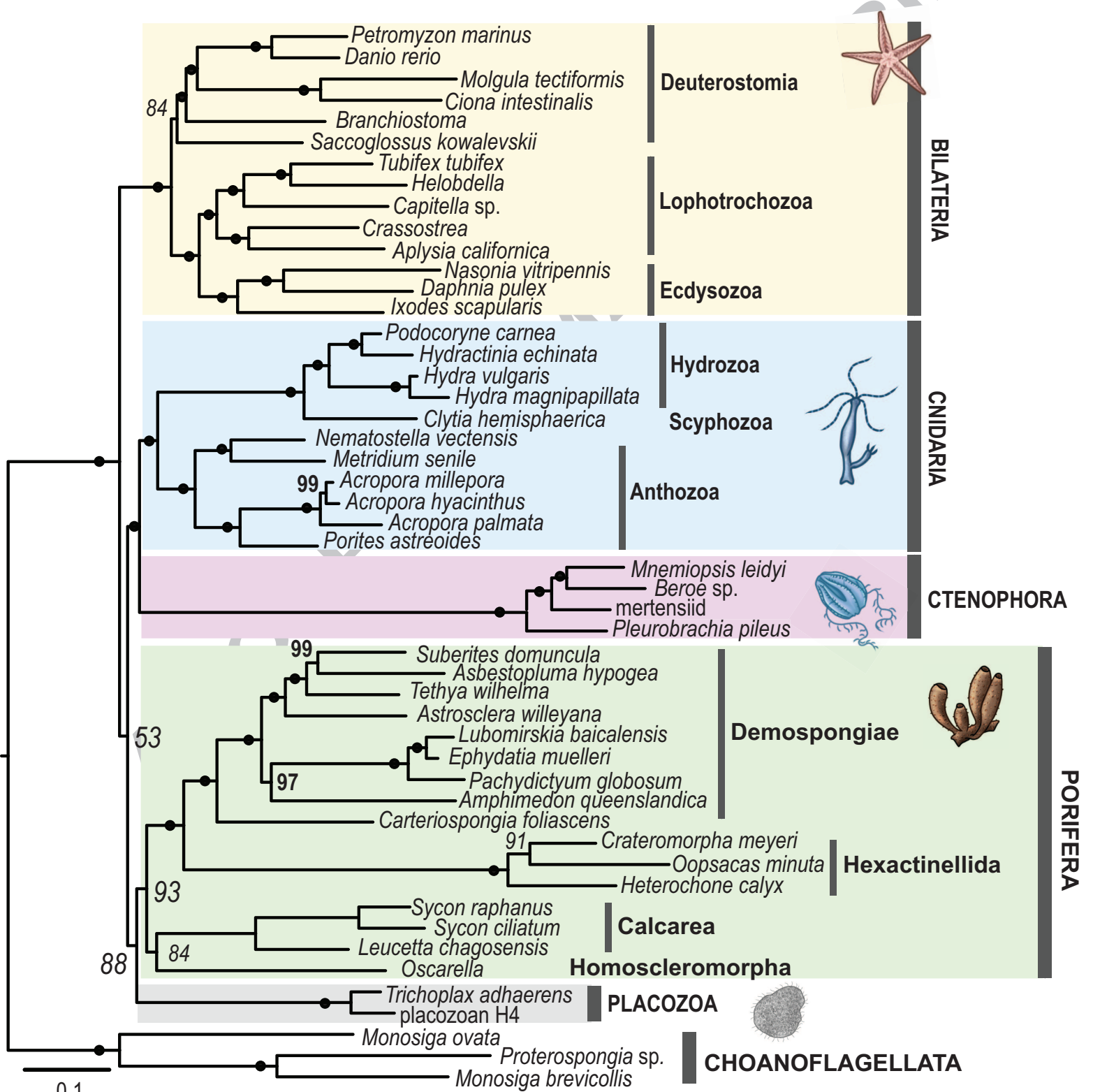
**Table 1. Large multi-gene matrices used for addressing the early metazoan phylogeny question**

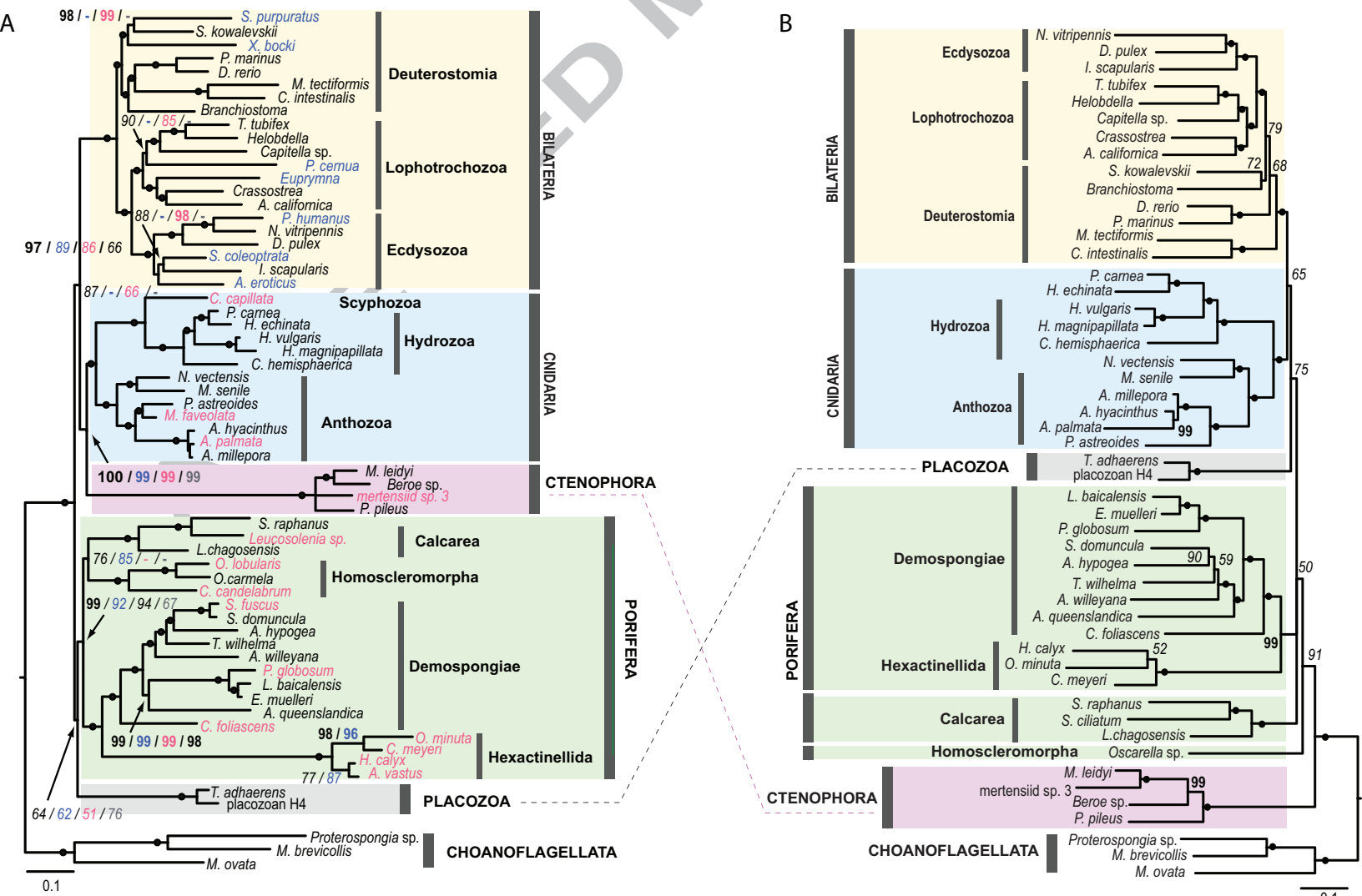
Gene matrix	Taxon #	Gene #	Matrix length (aa)	Variable site #	Allowed % missing data per taxon	Missing characters total (%)
Ribosomal <sup>a</sup>	63	87	14615	10445	95	28
	56	87	14615	10226	95	29
	49	87	14615	10288	50	14
	42	87	14615	10050	50	13
	50	78	11057	9538	80	16
Non-ribosomal <sup>a</sup>	63	35	9187	6322	95	36
	50	35	9187	6067	80	28
Combined 1 <sup>a</sup>	50	122	22975	15605	80	24
Combined 2 <sup>a</sup>	50	113	20244	13784	80	22
Dunn <i>et al.</i> 2008	77	150	21152	18085	93	50
Philippe <i>et al.</i> 2009	55	128	30257	20790	90	27

Multi-gene matrices used in this study are compared with two previously published large datasets (Dunn *et al.* 2008; Philippe *et al.* 2009).

<sup>a</sup>All parameters are indicated for matrices that include a single outgroup, Choanoflagellata

Figure 1





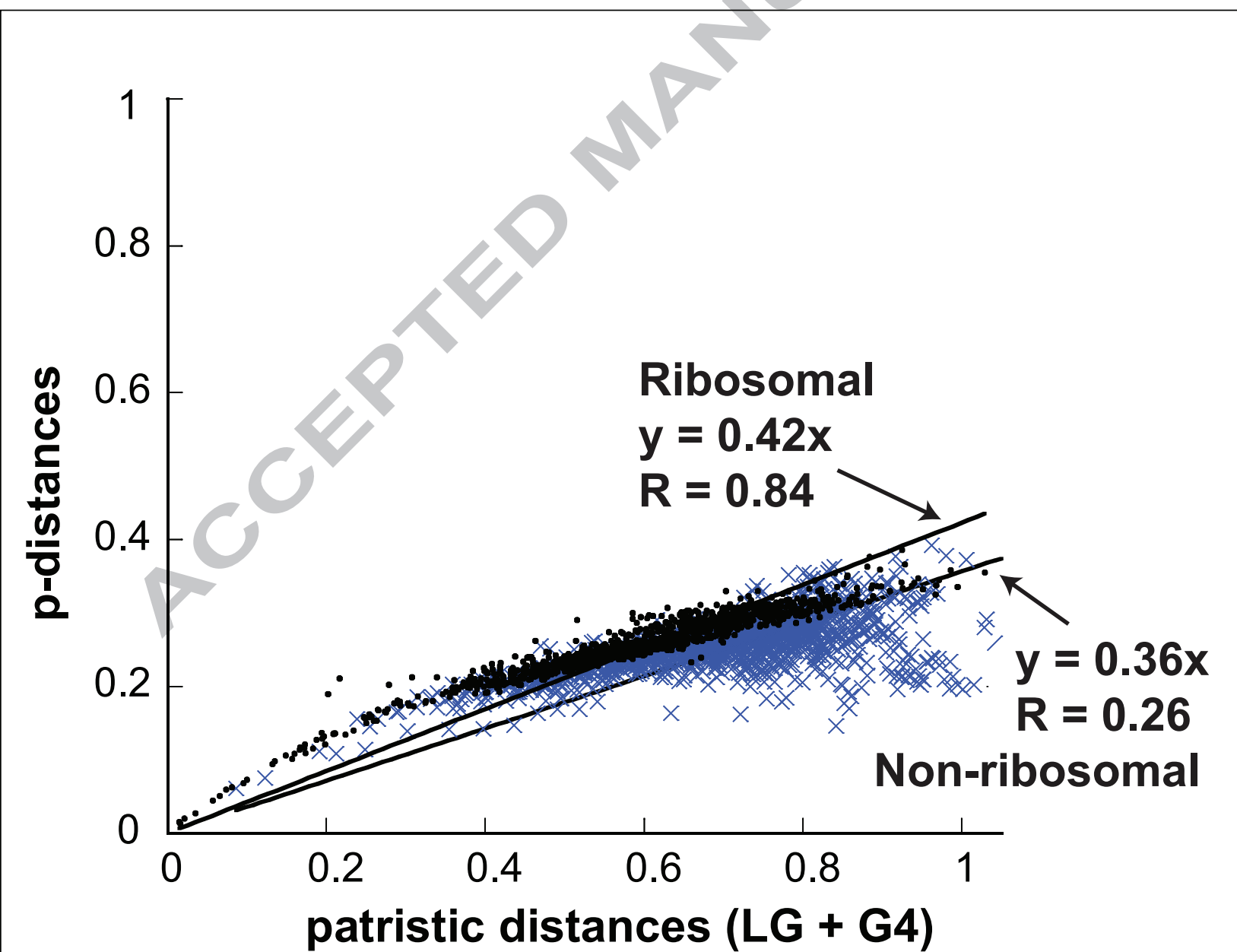
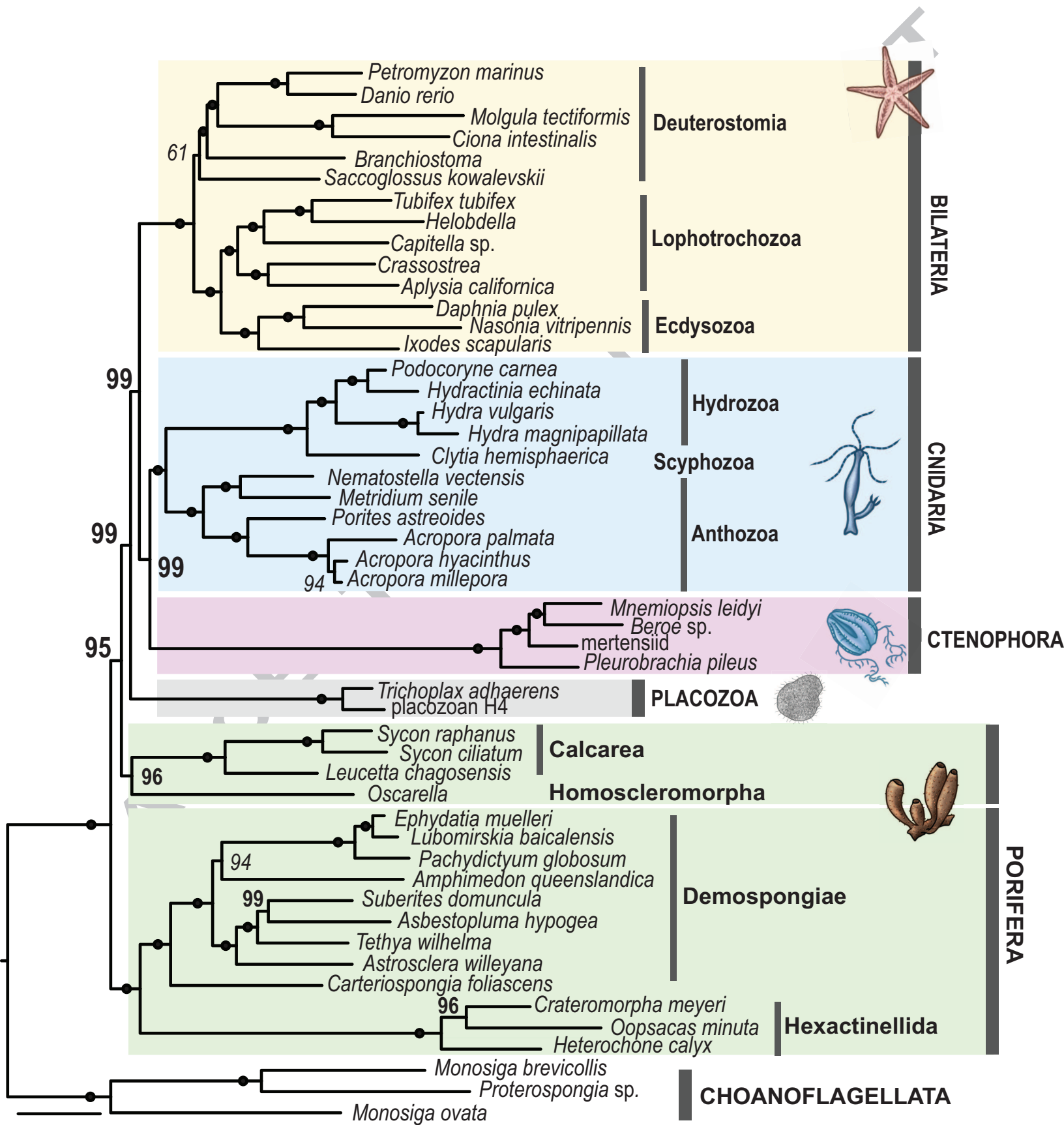
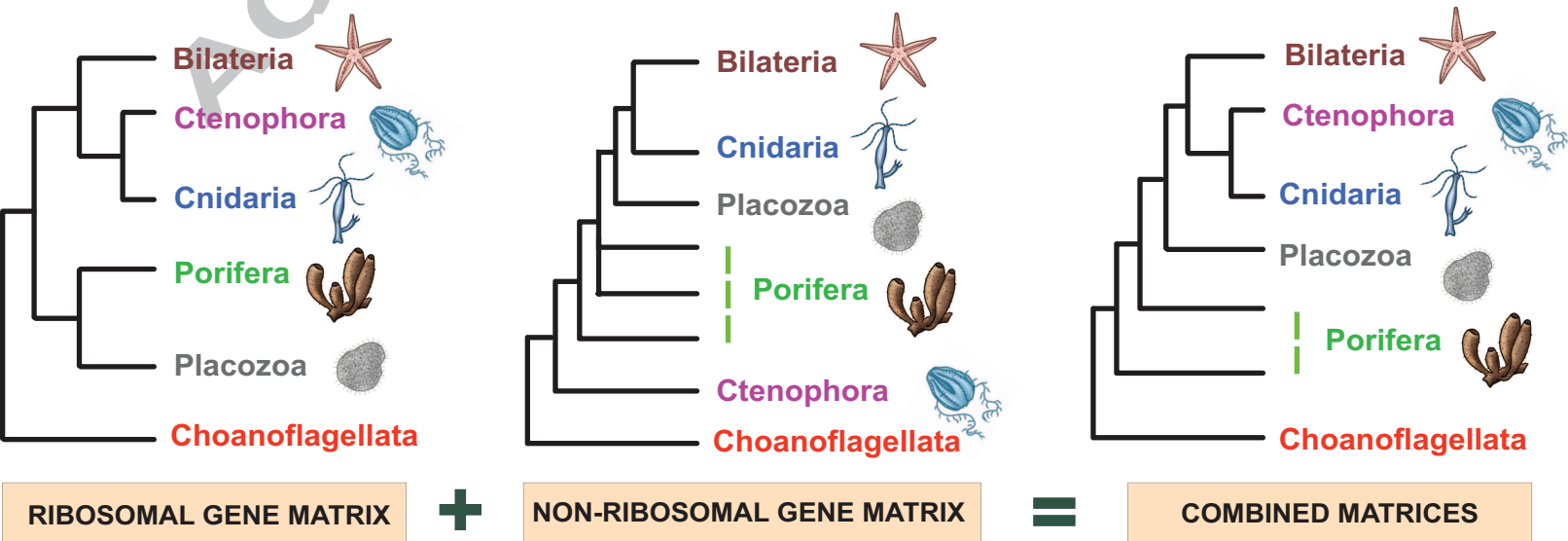


Figure 4







**Highlights**

- Deep metazoan phylogeny was tested using non-overlapping multi-gene matrices.
- Different partitions produce conflicting phylogenies.
- Level of saturation and LBA artifacts depend on gene sampling strategy.
- Ctenophora-basal and the sponge paraphyly correlate with higher saturation.
- Genes involved in translation support the Coelenterata and monophyly of Porifera.